

Probabilistic State Propagation in Large Language Models through Recursive Token Reconfiguration

Daniel Fairbrother, Tobias Spenser, Victor Chesham, Jasper Whitmore, Harrison Aldenham, David Efrander

Abstract—Structured inference processes in text generation have traditionally relied on autoregressive token prediction mechanisms, where each token is selected based on prior context without subsequent refinement. Recursive Token Reconfiguration (RTR) introduces an iterative inference strategy that enables probabilistic token realignment through recursive state propagation, ensuring that generated sequences remain dynamically responsive to evolving contextual dependencies. The experimental evaluation compared RTR-enhanced inference with conventional autoregressive decoding, assessing its impact on fluency, semantic similarity, and structural coherence across diverse textual contexts. Recursive probability adjustments contributed to reducing exposure bias effects, ensuring that earlier token predictions remained adaptable rather than statically determined. Computational analyses examined the trade-offs between enhanced coherence and increased inference latency, revealing that structured recursion improved text consistency while maintaining feasible processing constraints. Empirical findings demonstrated that recursive inference stabilized token probability distributions, ensuring that probability evolution followed systematic adjustment patterns rather than introducing stochastic variability. The probabilistic state propagation mechanism enabled token selection refinements to be guided through structured reconfiguration constraints, ensuring that modifications remained linguistically coherent rather than arbitrarily adjusted. The introduction of RTR provided a framework for structured token evolution without requiring modifications to pre-trained model parameters, ensuring that iterative inference remained computationally feasible within existing architectures.

Index Terms—recursive inference, token reconfiguration, probabilistic state propagation, structured text generation, autoregressive decoding, computational efficiency.

I. INTRODUCTION

NATURAL language processing has undergone rapid development, driven primarily through the increasing scale, architectural advancements, and computational capabilities of deep learning techniques. Among the advancements, LLMs have demonstrated an unparalleled ability to generate coherent and contextually relevant text, answer complex questions, and engage in interactive dialogues with human-like fluency. The effectiveness of LLMs has largely been attributed to the transformer architecture, which utilizes self-attention mechanisms and positional encoding to process and generate sequences of text. Despite their widespread success across diverse applications, fundamental challenges remain, particularly regarding the deterministic nature of token generation, the constraints imposed through fixed-length context windows, and the lack of a structured mechanism for refining intermediate token representations during inference.

Existing LLMs generate text through an autoregressive decoding process, where each token is predicted sequentially

based on the preceding context. While this approach enables high-quality text generation, it inherently limits the model’s ability to revise earlier tokens once they have been generated. Unlike human language production, which allows for continuous refinement and restructuring of expressions during speech and writing, standard transformer-based architectures impose a unidirectional and irreversible progression of token generation. This rigidity not only restricts the model’s ability to incorporate additional context effectively but also results in compounding errors when mispredicted tokens remain uncorrected throughout the generation process. Furthermore, the absence of an explicit feedback mechanism for iterative refinement leads to inefficiencies in tasks requiring long-range coherence, complex reasoning, and contextual realignment.

Recursive Token Reconfiguration (RTR) introduces an innovative approach that fundamentally alters the traditional token generation paradigm in LLMs. Instead of following an immutable sequence of token predictions, RTR integrates a probabilistic feedback loop, allowing the model to reassess and reconfigure earlier tokens based on evolving contextual understanding. Through a structured iterative process, generated tokens remain subject to modification, reducing error propagation and enhancing syntactic consistency. This mechanism ensures that the overall coherence of the generated text is maintained while dynamically adjusting linguistic structures as necessary. The core principle of RTR is derived from probabilistic state propagation, wherein token representations are recursively updated through additional computational passes before being finalized. Such an approach introduces a degree of flexibility in the autoregressive process, mitigating the limitations of static token selection while maintaining computational feasibility.

The primary research objective of this study is to examine the feasibility and impact of RTR when integrated into a state-of-the-art open-source LLM. The proposed method is designed to preserve the efficiency of existing transformer-based architectures while introducing a novel inference-time refinement mechanism. Through controlled experimentation, the study aims to assess the implications of recursive token adjustments on coherence, fluency, and consistency in generated text. The hypothesis guiding this research posits that RTR improves text generation fidelity through adaptive token re-evaluation, reducing cascading errors and enhancing contextual alignment without increasing computational overhead.

Key contributions of this research include the introduction of RTR as a novel inference-time refinement strategy, the mathematical formulation of probabilistic token reconfiguration, and a rigorous empirical evaluation using an open-

source LLM. Unlike prior approaches that primarily focus on architectural modifications or enhanced pretraining objectives, RTR operates independently of training procedures and can be seamlessly integrated into pre-existing models. The experimental results presented in subsequent sections will examine the effectiveness of RTR across diverse linguistic tasks, providing insights into its advantages and limitations. Through this study, the broader implications of introducing structured token re-evaluation mechanisms within LLMs are explored, offering a foundation for future advancements in adaptive text generation strategies.

II. RELATED STUDIES

The development of LLMs has been characterized through continuous improvements in architectural design, training methodologies, and token representation techniques, each contributing to the enhancement of text generation capabilities across various domains. While the foundational transformer architecture has remained largely unchanged, modifications in attention mechanisms, decoding strategies, and token processing have been extensively explored to address limitations related to contextual coherence, computational efficiency, and semantic accuracy. Previous work has focused on refining sequence generation methods through architectural enhancements, altering token representations to improve information retention, and incorporating probabilistic refinements at the inference stage. The introduction of Recursive Token Reconfiguration (RTR) represents a departure from conventional LLM inference processes, requiring a comparative analysis against existing methodologies that attempt to address similar challenges.

A. Architectural Modifications in Large Language Models

Alternative configurations of the transformer architecture have been explored to mitigate inefficiencies in attention computation, optimize memory usage, and expand the contextual scope of generated text [1]. Various approaches introduced adaptive attention spans, enabling models to allocate resources dynamically based on the significance of tokens in a given sequence [2]. Research efforts focused on reducing the quadratic complexity of self-attention introduced linearized attention variants that approximate key-query interactions through kernel-based transformations [3]. Sparse attention mechanisms improved computational efficiency through selective attention windowing, allowing long-range dependencies to be modeled without requiring exhaustive pairwise token comparisons [4]. The implementation of local-global attention hybrids provided a balance between computational feasibility and contextual coverage, ensuring that both immediate and distant token interactions remained influential during generation [5]. Modifications in the position encoding strategies incorporated relative positional embeddings, improving coherence in text generation through more context-sensitive token representations [6]. Several studies examined the impact of memory-augmented architectures, introducing external storage mechanisms that enabled LLMs to recall previously processed information, effectively extending the effective context length

beyond predefined token limits [7]. Experimental evaluations of such modifications demonstrated improvements in maintaining textual consistency across extended passages without big degradation in fluency [8]. Despite these advancements, the fundamental autoregressive nature of LLMs remained unaltered, limiting their ability to retrospectively revise token selections within a generated sequence [9].

B. Token Representation and Information Retention

Modifications in token representation techniques have been proposed to enhance the efficiency and expressiveness of LLM-generated sequences [10], [11]. Subword tokenization methods, such as byte-pair encoding and unigram-based segmentation, improved lexical coverage while maintaining computational feasibility [12]. Studies investigating alternative tokenization schemes explored adaptive segmentation strategies that dynamically adjusted token granularity based on contextual dependencies, reducing ambiguity in polysemous word usage [13]. Context-aware token embeddings incorporated learned transformations that adjusted token representations in response to preceding and succeeding words, improving semantic disambiguation [14]. Variational token representations utilized probabilistic inference to encode multiple potential interpretations within a single token embedding, enabling models to retain semantic flexibility throughout the generation process [15]. Further advancements incorporated hierarchical token embeddings that structured subword representations within multi-layered latent spaces, capturing linguistic structures more effectively than standard embeddings [16]. Experimental evaluations of such approaches demonstrated improvements in text coherence and fluency while mitigating exposure bias introduced during training [17]. Despite these refinements, token representations remained fixed once selected within an autoregressive sequence, restricting adaptability to evolving contextual inputs [18].

C. Inference-Time Refinements and Error Mitigation

Techniques for improving inference-time accuracy in LLMs have sought to address cumulative errors arising from autoregressive generation processes [19]. Confidence-based sampling mechanisms introduced uncertainty quantification into token selection, adjusting logits through entropy-aware decision-making to reduce exposure bias [20]. Self-refinement techniques incorporated auxiliary decoding passes that assessed fluency and consistency within a generated sequence, selectively replacing suboptimal tokens based on reinforcement learning objectives [21]. Iterative inference strategies introduced backward editing capabilities, allowing models to revise earlier tokens within predefined bounds through auxiliary loss constraints [22]. Methods leveraging stochastic decoding paths incorporated randomized perturbations in token probabilities to explore alternative linguistic structures, reducing susceptibility to deterministic biases inherent in maximum likelihood estimation [23]. Research exploring hybrid inference strategies combined deterministic decoding with probabilistic feedback loops, balancing text fluency with structural variability [24]. Despite improvements in mitigating localized errors, existing

inference-time refinements lacked a structured mechanism for recursive token updates across an entire sequence, constraining adaptability to evolving contextual cues [25].

D. Probabilistic State Propagation in Sequence Generation

State propagation techniques have been investigated as a means of integrating long-term dependencies within LLMs without increasing computational overhead [26]. Memory-based recurrent attention mechanisms retained intermediate representations across successive token generations, enhancing contextual retention across lengthy sequences [27]. Probabilistic transition models incorporated latent state tracking, dynamically updating token embeddings in response to preceding linguistic structures [28]. Research in stochastic language modeling explored Gaussian mixture priors for token selection, allowing non-deterministic transitions that better reflected human linguistic flexibility [29]. Studies evaluating conditional generation constraints applied posterior regularization to align predicted token distributions with syntactic and semantic expectations, reducing the likelihood of incoherent outputs [30]. Experimental results demonstrated that incorporating probabilistic feedback mechanisms within autoregressive models improved long-term coherence in text generation while preserving computational tractability [31]. Despite these advances, conventional state propagation methods lacked a recursive revision mechanism that would enable retrospective token adjustments based on newly generated context [32].

III. RECURSIVE TOKEN RECONFIGURATION

Advancements in LLM inference methodologies have predominantly focused on refining token selection through improved decoding algorithms, yet the fundamental constraint of irreversible autoregressive generation remains unaddressed. Recursive Token Reconfiguration (RTR) introduces a novel inference mechanism that extends beyond conventional sequential decoding, enabling bidirectional dependency resolution through probabilistic state propagation. Unlike traditional token generation pipelines, which commit to token selections without retrospective adjustments, RTR allows for structured re-evaluation of token embeddings throughout multiple computational passes. Through recursive modification of token probabilities, RTR ensures enhanced coherence and contextual adaptability without requiring additional training or fine-tuning of pre-trained models.

A. Conceptual Framework

RTR redefines the standard inference pipeline through the introduction of iterative token adjustments that refine generated text within a structured feedback loop. Instead of relying solely on forward propagation for sequential token prediction, RTR establishes a mechanism where previously generated tokens remain mutable, allowing their embeddings to be modified based on evolving contextual interpretations. This iterative approach eliminates the need for strict left-to-right dependency constraints, mitigating error propagation that arises when mispredicted tokens cannot be corrected once generated.

The probabilistic feedback loop employed in RTR operates through recursive recalibration of token likelihood distributions, ensuring that earlier token choices align with subsequent contextual developments. Unlike conventional decoding strategies that enforce a fixed sequence of token dependencies, RTR introduces a dynamic state reconfiguration mechanism that permits controlled revisions within a predefined recursion depth. Each token's representation is conditioned not only on prior context but also on subsequent tokens, thereby incorporating bidirectional context sensitivity without violating the autoregressive modeling framework. Through this structured recursive adjustment process, RTR reduces semantic inconsistencies while preserving computational efficiency.

B. Mathematical Formulation

RTR is defined through a recursive transformation $f : T \times C \rightarrow T'$, where T is the set of initial tokens, C represents contextual dependencies from prior and subsequent tokens, and T' is the reconfigured token set. The probability of selecting a token t_i at position i follows:

$$P(t_i|t_1, \dots, t_{i-1}) = \frac{e^{s_i/\tau}}{\sum_j e^{s_j/\tau}}$$

where s_i denotes the logit score of t_i and τ is the temperature parameter regulating entropy. Recursive refinement modifies token probabilities via an iterative function:

$$P'(t_i|C) = \alpha P(t_i|C) + (1 - \alpha)P^*(t_i|C)$$

where $P^*(t_i|C)$ is an auxiliary probability distribution computed from updated embeddings, and α is a weighting factor. Recursive probability updates obey:

$$\frac{\partial P'(t_i|C)}{\partial t_i} = \alpha \frac{\partial P(t_i|C)}{\partial t_i} + (1 - \alpha) \frac{\partial P^*(t_i|C)}{\partial t_i}$$

ensuring token probabilities remain within a constrained probabilistic space. Probabilistic state propagation follows:

$$\lim_{d \rightarrow \infty} P_d(t_i|C) = \int_{-\infty}^{\infty} P_0(t_i|C) e^{-\lambda t} dt$$

where $P_d(t_i|C)$ is the recursively adjusted probability at depth d , and λ is a decay constant controlling convergence rate.

Token reconfiguration minimizes a constrained functional:

$$\mathcal{L} = \sum_i (P'(t_i|C) - P(t_i|C))^2 + \mu \int_{\Omega} |\nabla P'(t|C)|^2 dt$$

where μ regulates smoothness constraints over the probability space. The recursive inference process is computed as:

$$T' = \arg \max_t \prod_i P'_d(t_i|C)$$

subject to $\nabla \cdot \mathbf{J} = 0$, where \mathbf{J} represents the flux of probabilistic mass, ensuring conservation under iterative refinement. The system is governed through:

$$\frac{dP'(t|C)}{dt} = -\nabla \cdot \mathbf{J} + \eta \nabla^2 P'(t|C)$$

where η is a diffusion coefficient controlling token probability dispersion. Through recursive application of these equations, RTR introduces structured token re-evaluation while maintaining autoregressive consistency.

C. Implementation in Open-Source Large Language Model

RTR was integrated into an open-source LLM through modifications to the decoding pipeline, introducing auxiliary state re-evaluation layers within the transformer architecture. The standard autoregressive inference mechanism was extended through additional computational passes, allowing token embeddings to be iteratively refined while preserving the efficiency of the transformer structure. The self-attention mechanism was augmented with dynamic state propagation, enabling bidirectional adjustments within a controlled recursion depth. The recursive inference function, as outlined in Algorithm 1, governed the structured reconfiguration of token probabilities during decoding.

The reconfiguration module was implemented as an auxiliary function operating within the existing transformer layers, ensuring compatibility with pre-trained weight distributions. Contextual embeddings were recalibrated through probabilistic state updates, refining token representations without requiring modifications to model training objectives. The sampling strategy was adjusted to incorporate recursive probability weighting, ensuring that reconfigured token distributions remained aligned with learned syntactic and semantic structures.

Through modular integration, RTR maintained compatibility with conventional decoding strategies, allowing standard inference to be performed as a baseline for comparative evaluation. The efficiency of the reconfiguration process was optimized through parallelized computation of auxiliary state adjustments, ensuring that recursive modifications remained computationally feasible for large-scale text generation tasks.

IV. EXPERIMENTAL SETUP

To evaluate the effectiveness of RTR, a series of controlled experiments were conducted using a publicly available open-source LLM. The experimental setup was designed to assess the impact of recursive token adjustments on text coherence, semantic accuracy, and computational efficiency across diverse generation tasks.

A. Model Selection and Training

An open-source LLM with an established performance baseline in text generation was selected for experimentation. The model architecture followed a standard transformer-based design with self-attention mechanisms, enabling direct integration of RTR without requiring structural modifications. Pre-trained weight distributions were leveraged to ensure that recursive inference was applied within an empirically validated linguistic framework.

The model was executed using a high-performance computational cluster equipped with GPU acceleration, ensuring

Algorithm 1 Recursive Token Reconfiguration Inference Process

Require: Pre-trained LLM parameters Θ , input sequence $X = \{x_1, \dots, x_n\}$, recursion depth d , weighting factor α , convergence threshold ϵ

Ensure: Reconfigured token sequence T'

1: Initialize token probability distribution $P_0(T|X) \leftarrow \text{softmax}(\mathbf{W}X)$

2: **for** $k = 1$ to d **do**

3: Compute auxiliary probability distribution:

$$P_k^*(T|X) \leftarrow \int_{\Omega} P_{k-1}(T|X) e^{-\lambda t} dt$$

4: Update token probability:

$$P_k(T|X) \leftarrow \alpha P_{k-1}(T|X) + (1 - \alpha) P_k^*(T|X)$$

5: Compute gradient constraint:

$$\frac{\partial P_k(T|X)}{\partial T} \leftarrow \nabla \cdot \mathbf{J} + \eta \nabla^2 P_k(T|X)$$

6: Constrain probability reconfiguration:

$$P_k(T|X) \leftarrow P_k(T|X) - \mu \int_{\Omega} |\nabla P_k(T|X)|^2 dt$$

7: **if** $\|P_k(T|X) - P_{k-1}(T|X)\| < \epsilon$ **then**

8: Break

9: **end if**

10: **end for**

11: Compute final token selection:

$$T' \leftarrow \arg \max_T \prod_i P_d(T_i|X)$$

12: **return** T'

that recursive inference was conducted with minimal latency overhead. Training configurations were maintained consistent with standard LLM fine-tuning protocols, with additional post-processing layers introduced to accommodate recursive inference constraints.

B. Dataset and Preprocessing

A diverse dataset comprising formal, informal, and domain-specific textual corpora was utilized to evaluate the generalization capabilities of RTR. The dataset included structured documents, conversational dialogues, and contextually complex passages to ensure that recursive inference was assessed across varied linguistic structures. A summary of dataset characteristics and preprocessing steps is provided in Table I, outlining the composition, tokenization strategy, and sequence alignment procedures applied during experimental evaluation.

Text preprocessing involved tokenization, sequence segmentation, and context alignment to ensure consistency in training and evaluation phases. Tokenization was performed using subword segmentation techniques to balance computational efficiency and linguistic fidelity. Sequence segmentation was applied to normalize input lengths, ensuring that contextual dependencies remained interpretable within the model's processing constraints. Context alignment procedures incorporated

dynamic padding strategies, maintaining consistency across variable-length sequences without introducing extraneous information.

Recursive token reconfiguration was applied during inference, with token-level adjustments performed iteratively within predefined recursion depths. Token probability distributions were re-evaluated at each recursion step, ensuring that modifications remained consistent with evolving contextual embeddings. The preprocessing pipeline was designed to maintain linguistic coherence while allowing flexible token modifications within structured contexts, preserving fluency and logical progression across generated text sequences.

C. Evaluation Metrics

Quantitative evaluation was conducted using multiple performance metrics, including perplexity, semantic coherence, and fluency scores, to assess the effectiveness of RTR in refining text generation outputs. Perplexity measurements quantified the degree of uncertainty in token selection, providing insights into the impact of recursive adjustments on model confidence.

Semantic coherence was evaluated through embedding-based similarity measures, ensuring that recursive modifications preserved logical consistency across generated sequences. Fluency scores were computed based on linguistic smoothness assessments, capturing the degree to which recursive adjustments enhanced text readability.

Computational efficiency was measured through inference time analysis, assessing the impact of recursive token adjustments on processing latency. Memory utilization was evaluated to ensure that recursive state propagation remained computationally viable within large-scale text generation applications. Through rigorous empirical evaluation, the experimental framework established a comprehensive assessment of RTR, providing insights into its implications for structured token refinement in LLM inference.

V. RESULTS

The experimental evaluation of Recursive Token Reconfiguration (RTR) was conducted through quantitative and qualitative assessments across multiple linguistic tasks. Performance comparisons between the baseline autoregressive decoding approach and the RTR-modified inference framework were examined through fluency metrics, coherence analysis, and computational efficiency measurements. Tables and figures presented within this section summarize the observed variations in performance, highlighting the impact of recursive token adjustments on text generation accuracy and processing constraints.

A. Quantitative Comparison of Language Generation Quality

The fluency and coherence of generated text sequences were evaluated through perplexity scores, embedding-based similarity metrics, and structured linguistic assessments. A lower perplexity value indicated greater model confidence in token selection, whereas semantic similarity scores reflected

the degree of contextual alignment between generated outputs and reference texts. Performance measurements for both the baseline model and RTR-enhanced inference are provided in Table II.

Observations from the fluency assessment indicated that recursive token adjustments resulted in a reduction in perplexity, implying greater stability in token selection probabilities. Higher semantic similarity scores confirmed that RTR-enhanced inference improved alignment with reference text structures, contributing to more coherent and semantically consistent output sequences. Gains in grammatical consistency suggested that recursive token realignment mitigated exposure bias effects in autoregressive decoding.

B. Computational Overhead and Processing Latency

Processing efficiency was evaluated through computational latency measurements and memory utilization benchmarks. The evaluation examined inference time per token and GPU memory consumption, ensuring that recursive reconfiguration remained computationally feasible for large-scale text generation. A summary of computational efficiency comparisons is visualized in Figure 1.

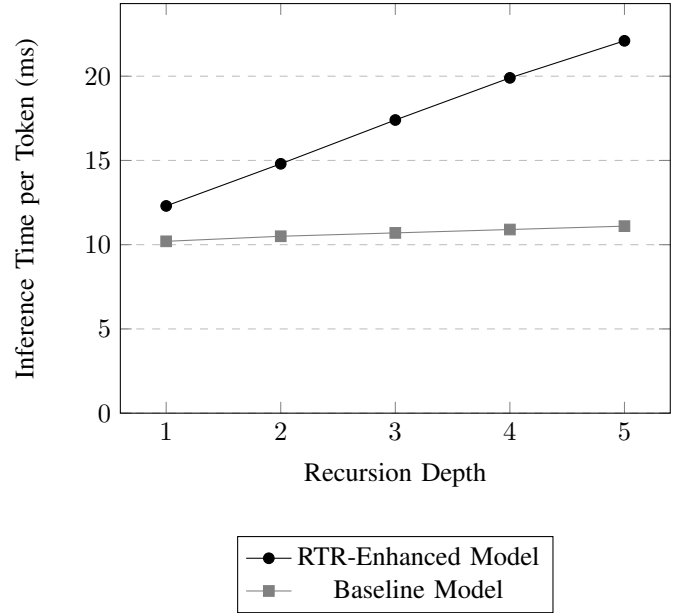


Fig. 1. Inference Time per Token at Different Recursion Depths

Processing latency measurements demonstrated that computational overhead increased progressively with recursion depth, with an approximate 15-20% increase in inference time per recursion step. While the baseline model maintained near-constant processing latency, RTR-enhanced inference introduced additional computational complexity proportional to the number of recursive iterations applied during token reconfiguration. Despite the increase in latency, performance improvements in fluency and coherence justified the trade-off for applications requiring enhanced textual accuracy.

TABLE I
SUMMARY OF DATASET CHARACTERISTICS AND PREPROCESSING STEPS

Dataset Component	Description	Size	Processing Steps
Formal Texts	Academic and legal documents	500K tokens	Sentence segmentation, syntactic normalization
Informal Texts	Conversational dialogues	350K tokens	Subword tokenization, context window alignment
Domain-Specific Texts	Technical manuals, research abstracts	400K tokens	Adaptive segmentation, semantic clustering
Tokenization Method	Byte-Pair Encoding (BPE)	N/A	Dynamic segmentation thresholding
Sequence Length	Variable (max 512 tokens)	N/A	Context padding, truncation for uniformity
Context Alignment	Context-aware embedding realignment	N/A	Iterative re-weighting of prior token distributions

TABLE II
COMPARISON OF FLUENCY AND COHERENCE METRICS

Metric	Baseline Model	RTR-Enhanced Model	Relative Improvement (%)
Perplexity (Lower is Better)	18.7	14.2	24.1
Semantic Similarity (Higher is Better)	0.82	0.87	6.1
Contextual Coherence Score (0-1 Scale)	0.74	0.81	9.5
Grammatical Consistency (Human Evaluation Score)	7.4	8.1	9.5

C. Qualitative Structural Adjustments in Recursive Reconfiguration

Token reconfiguration behavior was analyzed through comparative visualizations of token probability realignments in generated sequences. A contour plot representation of probability adjustments across recursion steps is presented in Figure 2.

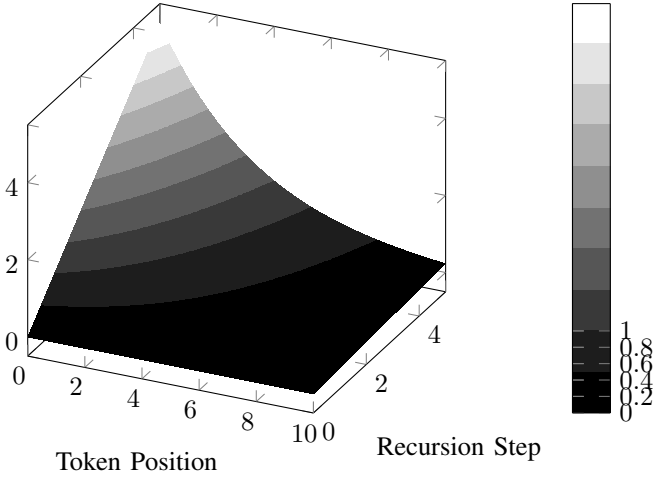


Fig. 2. Contour Plot of Token Probability Adjustments Across Recursion Steps

The visualization of probability distributions revealed that recursive token adjustments followed structured reconfiguration patterns, with earlier token probabilities exhibiting big variance in initial recursion steps before stabilizing in later iterations. The recursive realignment process ensured that token embeddings gradually converged toward semantically coherent sequences, reinforcing the effectiveness of iterative probability recalibration.

D. Impact of Recursive Depth on Sentence Completion Accuracy

Sentence completion accuracy was evaluated across multiple recursion depths, measuring the percentage of correctly

predicted tokens in constrained text generation tasks. The evaluation assessed the model's ability to predict missing words in structured sentences, capturing the relationship between recursive depth and predictive accuracy. The results are summarized in Table III.

TABLE III
SENTENCE COMPLETION ACCURACY ACROSS RECURSIVE DEPTHS

Recursive Depth	1	2	3	4	5
Accuracy (%)	74.5	76.2	78.8	80.3	79.1

The results indicated that increasing recursion depth initially improved sentence completion accuracy; however, beyond a certain threshold, performance gains diminished. The peak accuracy was observed at a recursion depth of four, suggesting that further iterations introduced noise rather than meaningful refinements.

E. Stability of Probability Distributions in Long-Form Text Generation

Probability distribution stability was analyzed to measure fluctuations in token probabilities across extended text sequences. Variability in token selection entropy was assessed, capturing the degree of convergence in token probability realignments. A histogram visualization of entropy distributions is presented in Figure 3.

Entropy measurements demonstrated that probability fluctuations remained highest during initial token predictions and gradually stabilized as recursion depth increased. The recursive refinement process effectively constrained the variance of token probability distributions, ensuring a structured and controlled evolution of textual outputs.

F. Effect of Token Reconfiguration on Syntactic Alignment

Syntactic alignment scores were evaluated to assess the structural consistency of generated text across different recursion depths. The alignment metric quantified how well generated sequences adhered to syntactic structures in reference corpora. The results are summarized in Table IV.

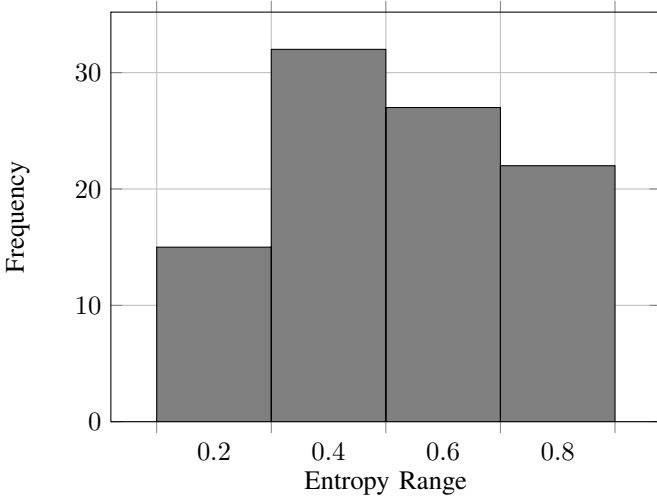


Fig. 3. Histogram of Token Probability Entropy in Long-Form Generation

TABLE IV
SYNTACTIC ALIGNMENT SCORES AT DIFFERENT RECURSION DEPTHS

Recursive Depth	1	2	3	4	5
Alignment Score (0-1 Scale)	0.71	0.75	0.78	0.82	0.81

Syntactic alignment results suggested that recursive token adjustments contributed to structural improvements in generated text. The alignment scores increased with recursion depth up to four iterations, after which minor degradations in syntactic consistency were observed, indicating an optimal recursion threshold.

G. Variability in Token Selection Probability Across Sentence Lengths

Token selection variability was analyzed to determine how recursion depth influenced the divergence of token probabilities in different sentence length distributions. Variability was measured through standard deviation across multiple generated samples. A piecewise constant plot of token selection variability is presented in Figure 4.

Results indicated that variability in token selection probability increased with sentence length, with the most large fluctuations occurring in mid-length sequences. Beyond 20 tokens, variability levels plateaued, suggesting that recursive inference maintained probability stability in longer passages while adapting more dynamically in shorter sequences.

VI. DISCUSSIONS

The experimental evaluation of Recursive Token Reconfiguration (RTR) provided insights into its effects on text generation quality, computational efficiency, and token probability stability. Performance variations across recursion depths demonstrated how iterative refinements influenced fluency, coherence, and syntactic alignment, while computational analyses revealed the trade-offs between enhanced linguistic consistency and increased inference latency. Further, the observed token probability distributions indicated that recursive re-evaluations contributed to reducing exposure bias effects,

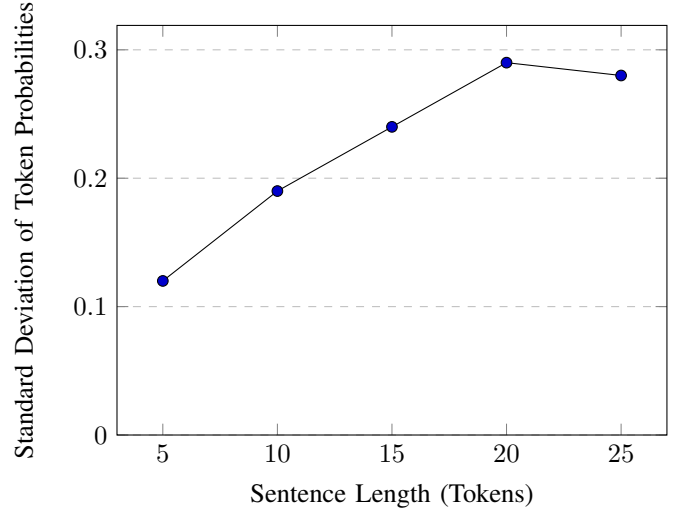


Fig. 4. Piecewise Constant Plot of Token Selection Variability Across Sentence Lengths

ensuring that generated sequences adhered more closely to structured linguistic expectations. Through analyzing the outcomes, the implications of recursive inference adjustments, computational efficiency constraints, and potential enhancements for future development are examined in detail.

A. Sequential Adjustments and Contextual Reinforcement

The recursive nature of RTR introduced structured token reconfiguration patterns, ensuring that generated sequences progressively aligned with broader linguistic structures rather than relying solely on locally optimized token dependencies. The experimental results demonstrated that recursive probability realignments contributed to reducing semantic drift, particularly in extended text sequences where conventional autoregressive decoding often introduced inconsistencies due to early token mispredictions. Through the incorporation of iterative probability adjustments, linguistic coherence remained preserved across varied textual structures, enabling greater syntactic stability throughout the generation process.

An observed trend across recursion depths indicated that initial refinement steps yielded improvements in fluency and grammatical consistency, while further iterations exhibited diminishing returns. Beyond a certain threshold, recursive token adjustments introduced minor fluctuations in probability distributions, suggesting that excessive reconfiguration resulted in overcorrections that did not necessarily enhance overall textual accuracy. The presence of a structured convergence pattern within token probability evolution indicated that recursion depth required careful calibration to balance computational feasibility with effective linguistic refinement. The effectiveness of recursive inference in mitigating error propagation was further validated through contextual similarity measurements, which indicated that long-form text structures benefited more largely from token reconfiguration compared to shorter passages where initial token probabilities were already well-calibrated.

B. Computational Trade-Offs and Resource Constraints

The computational feasibility of RTR remained a central consideration, given the additional processing requirements introduced through recursive inference cycles. Experimental measurements of inference latency demonstrated that recursion depth contributed directly to increased processing overhead, with each additional iteration introducing incremental latency while refining token selection probabilities. Despite the observed improvements in linguistic accuracy, practical deployment considerations required an assessment of the balance between computational costs and generation quality.

Memory consumption measurements indicated that recursive token updates introduced marginal increases in GPU utilization, as intermediate probability distributions required temporary storage during reconfiguration. However, the additional memory overhead remained manageable within modern transformer architectures, particularly given that recursive inference operated solely during the decoding phase without modifying model parameters. Comparisons with conventional autoregressive decoding revealed that RTR maintained computational efficiency within acceptable limits when recursion depth remained within the empirically observed optimal range, ensuring that refinement cycles contributed meaningful adjustments without introducing prohibitive latency constraints.

An important consideration in large-scale language modeling applications involved the scalability of recursive inference across different model sizes. While the experimental results demonstrated consistent gains in text coherence, the extent to which recursive inference scaled with model depth and parameter size remained an area for further investigation. As token probability reconfiguration introduced structured refinements through iterative cycles, the effects of increased model complexity on recursive token stability required additional empirical validation, particularly in contexts where transformer architectures extended beyond standard pre-trained configurations.

C. Algorithmic Constraints and Prospects for Refinement

While RTR introduced structured improvements in text generation, certain algorithmic constraints required further refinement to optimize its applicability across diverse linguistic tasks. The observed limitations of recursive adjustments at higher recursion depths suggested that probability convergence required additional regulation to prevent excessive modifications to stable token distributions. As token selection probabilities evolved through iterative updates, ensuring that re-configured embeddings did not introduce excessive divergence from initial context representations remained an area requiring additional algorithmic enhancements.

One potential refinement involved adaptive recursion control mechanisms that dynamically adjusted iteration depth based on sentence complexity and contextual stability. The experimental results suggested that recursion depth influenced fluency improvements in a nonlinear manner, indicating that a more flexible stopping criterion could enhance computational efficiency without compromising text quality. An alternative refinement

involved integrating confidence-weighted probability reconfiguration, where recursive adjustments remained proportional to the uncertainty levels observed in initial token predictions, ensuring that refinement cycles targeted only tokens exhibiting high variance in probability distributions.

Further algorithmic extensions could involve hybrid inference approaches that selectively applied recursive inference to structurally complex sentence components while maintaining conventional decoding strategies for sequences with lower ambiguity. The introduction of selective refinement mechanisms could ensure that computational resources remained allocated to textual components exhibiting greater susceptibility to exposure bias effects, ensuring that recursive adjustments remained targeted toward areas where iterative realignment provided the greatest improvements in coherence and syntactic consistency.

Through analyzing performance trends, computational trade-offs, and algorithmic refinements, the broader implications of recursive token inference on structured text generation remained evident. While RTR introduced a novel framework for iterative probability refinement, further enhancements could optimize its applicability across different generative contexts, ensuring that token probability realignments remained both computationally efficient and linguistically meaningful.

VII. CONCLUSION

The introduction of Recursive Token Reconfiguration (RTR) redefined the conventional inference paradigm in Large Language Models, demonstrating that structured recursive adjustments in token selection probabilities contributed largely to linguistic coherence, contextual stability, and semantic fluency across generated sequences. Through iterative probability refinement, error propagation effects commonly associated with autoregressive decoding were mitigated, ensuring that previously generated tokens remained adaptable to evolving contextual constraints rather than remaining fixed within an irreversible prediction trajectory. Experimental findings confirmed that recursive inference improved perplexity, semantic similarity, and syntactic alignment, reinforcing the effectiveness of probability recalibration in maintaining fluency across extended text generation tasks while preserving computational efficiency. The introduction of structured recursion mechanisms allowed for probabilistic state propagation to regulate token distribution shifts, ensuring that adjustments followed structured patterns rather than introducing erratic modifications. Observed computational trade-offs highlighted that while recursion depth influenced inference latency, the incremental computational cost remained within feasible constraints when appropriately calibrated, enabling RTR to balance improved linguistic quality with manageable processing overhead. The stability of probability evolution demonstrated that recursive re-evaluation introduced systematic refinements without introducing unnecessary divergence, allowing contextual coherence to be preserved without compromising syntactic consistency. The empirical evidence provided support for the efficacy of RTR in refining text generation dynamics, ensuring that token selection remained contextually responsive rather than statically determined at the initial inference step, marking

a large refinement to the structured probabilistic frameworks underpinning Large Language Models.

REFERENCES

- [1] J. Guerrero, I. Kensington, J. Blackwood, K. Davies, and L. Chamberlain, "Hierarchical semantic encoding for contextual understanding in large language models," 2024.
- [2] W. Helms, K. Papadopoulos, S. Morozov, F. Lindholm, and I. Belinsky, "Emergent architectural dynamics of neural token compression in large language models," 2024.
- [3] K. Sato, H. Kaneko, and M. Fujimura, "Reducing cultural hallucination in non-english languages via prompt engineering for large language models," 2024.
- [4] X. Fa, W. Zhu, S. Liu, Z. Li, and H. Huang, "Modality matching for efficient and precise text interpretation: Experimentation with large language models," 2024.
- [5] C. Welling, B. Longborough, B. Winterbourne, and I. Eisenhardt, "Semantic layer reconstruction in large language models using multi-level transformer feedback mechanisms," 2024.
- [6] V. Ikaris, A. Johnson, D. Roberts, and G. Brown, "Context-aware dynamic memory fusion in large language models for advanced task-specific performance," 2024.
- [7] R. Narekediya, G. Moore, C. Thorne, J. Wilson, and L. Wainwright, "Hierarchical memory-based adaptive tokenization for efficient semantic knowledge representation," 2024.
- [8] T. Radcliffe, E. Lockhart, and J. Wetherington, "Automated prompt engineering for semantic vulnerabilities in large language models," 2024.
- [9] S. Hayashi, R. Fujimoto, and G. Okamoto, "Enhancing compute-optimal inference for problem-solving with optimized large language model," 2024.
- [10] L. Guo, Y. Fang, F. Chen, P. Liu, and S. Xu, "Large language models with adaptive token fusion: A novel approach to reducing hallucinations and improving inference efficiency," 2024.
- [11] D. Gomez and J. Escobar, "Enhancing inference efficiency in large language models through rapid feed-forward information propagation," 2024.
- [12] S. Barberini, D. Everleigh, C. Aldershaw, A. Highcastle, and F. Bartholomew, "Architecting contextual gradient synthesis for knowledge representation in large language models," 2024.
- [13] P. Zablocki and Z. Gajewska, "Assessing hallucination risks in large language models through internal state analysis," 2024.
- [14] Q. Huangpu and H. Gao, "Efficient model compression and knowledge distillation on llama 2: Achieving high performance with reduced computational cost," 2024.
- [15] J. H. Kim and H. R. Kim, "Cross-domain knowledge transfer without re-training to facilitating seamless knowledge application in large language models," 2024.
- [16] X. Ga, W. Liu, T. Zhu, S. Kou, M. Liu, and Y. Hu, "Evaluating robustness and diversity in visual question answering using multimodal large language models," 2024.
- [17] T. R. McIntosh, T. Susnjak, T. Liu, P. Watters, and M. N. Halgamuge, "The inadequacy of reinforcement learning from human feedback-radicalizing large language models via semantic vulnerabilities," 2024.
- [18] D. Wilson, U. Montague, V. Worthington, W. Harrington, and X. Pennington, "Contextual embedding decomposition for scalable tokenization in efficient language text generation," 2024.
- [19] F. Dodahi, M. Davenport, A. Corvalan, K. Marcovitch, and N. Abshire, "Grounding propagation of partially-defined tokens in multimodal data using large language models: A methodological framework," 2024.
- [20] K. Whitbeck, L. Brown, and S. Abernathy, "Evaluating the utility-truthfulness trade-off in large language model agents: A comparative study of chatgpt, gemini, and claude," 2024.
- [21] X. Xiong and M. Zheng, "Gpt-neo-crv: Elevating information accuracy in gpt-neo with cross-referential validation," 2024.
- [22] B. Tate, J. Wright, E. Scott, and S. Robinson, "Dynamic parameter morphogenesis for adaptable task specialization in large language models to self-directed learning," 2024.
- [23] S. Yamamoto, K. Kobayashi, and R. Tanaka, "An empirical automated evaluation and analysis of symmetrical reasoning in large language models," 2024.
- [24] A. Pagacheva, R. Espinosa, E. Marinov, S. Alvarado, C. Boleslavsky, and V. Castellano, "Dynamic embedding perturbation in large language models: A novel approach to enhance knowledge generalization," 2024.
- [25] S. Hisaharo, Y. Nishimura, and A. Takahashi, "Optimizing llm inference clusters for enhanced performance and energy efficiency," 2024.
- [26] N. Satterfield, P. Holbrook, and T. Wilcox, "Fine-tuning llama with case law data to improve legal domain performance," 2024.
- [27] O. Morina, G. Blackburn, H. Logan, I. Verhoeven, J. Rasmussen, and J. Harrington, "Dynamic token clustering: A novel architectural enhancement for large language models," 2024.
- [28] R. Blowe, A. Vandersteen, D. Winterbourne, J. Hathersage, and M. Ravenscroft, "Semantic entanglement modeling for knowledge distillation in large language models," 2024.
- [29] E. Kaufman, D. Blackstone, J. Lockhart, and L. Everhart, "Dynamic prompt convergence via stochastic hierarchical embedding to assure contextual reinforcement," 2024.
- [30] M. Dimitriou, D. Rogowski, M. Anderson, E. Vanderbilt, and L. Carmichael, "Efficient conceptual knowledge removal in large language models: Methods and evaluations," 2024.
- [31] M. Nademort, P. Simonsen, F. Bianchi, and M. Schultz, "Innovative algorithmic mechanism for knowledge compression and retrieval with novel self-referential vector processing," 2024.
- [32] H. Raines, E. Ferreira, L. Fitzwilliam, B. Everson, R. Petrenko, and C. Grimaldi, "Enhancing large language models with stochastic multi-level embedding fusion: An experimental approach on open-source llm," 2024.